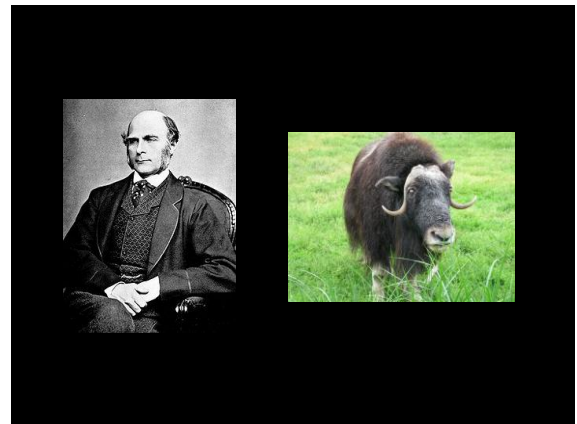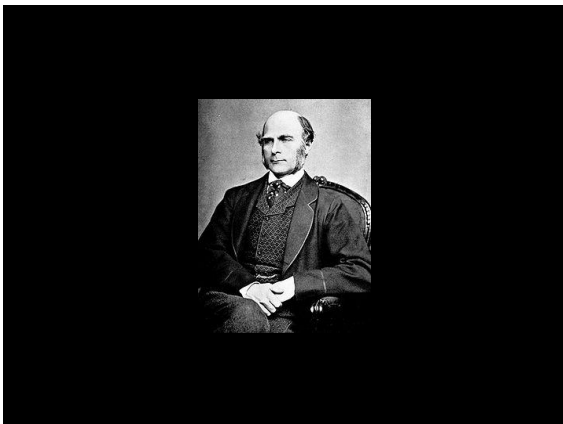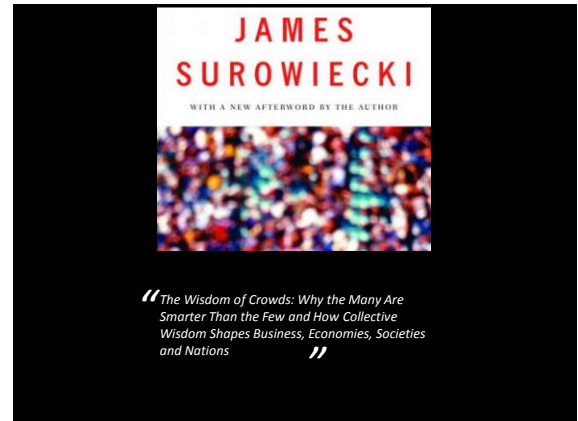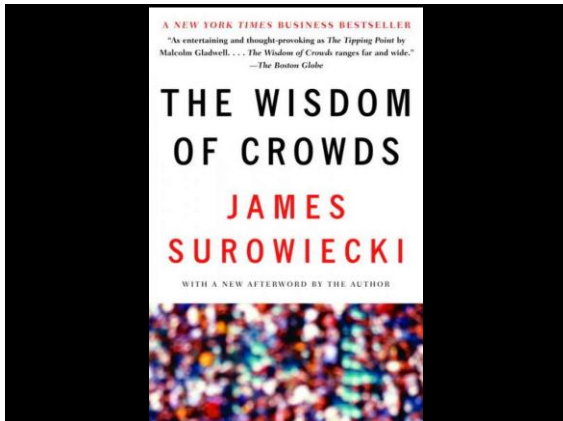# Ensemble Learning

CS4780/5780 – Machine Learning
Fall 2013

Igor Labutov
Cornell University

---

# Ensemble Learning

A class of "meta" learning algorithms

---



---



"The Wisdom of Crowds: Why the Many Are Smarter Than the Few and How Collective Wisdom Shapes Business, Economies, Societies and Nations"

---



---

## Slide 1

**1,198 lb**      **1,197 lb**



## Slide 2



| Criteria | Description |
|---|---|
| Diversity of opinion | *Each person should have private information even if it's just an eccentric interpretation of the known facts.* |
| Independence | *People's opinions aren't determined by the opinions of those around them.* |
| Decentralization | *People are able to specialize and draw on local knowledge.* |
| Aggregation | *Some mechanism exists for turning private judgments into a collective decision.* |

## Ensemble Learning

A class of "meta" learning algorithms

Combining multiple classifiers to increase performance

Very effective in practice

Good theoretical guarantees

Easy to implement!

## Ensemble

Problem : given *T* binary classification hypotheses ($h_1$,…, $h_T$), **find** a combined classifier:
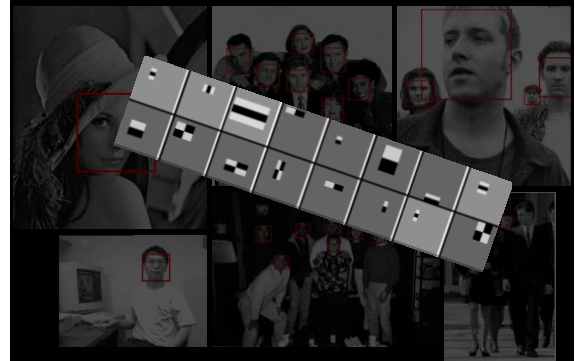
$$h_S(x) = \mathrm{sign}\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right)$$

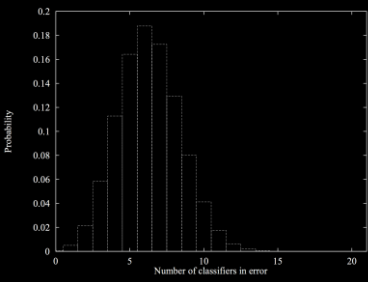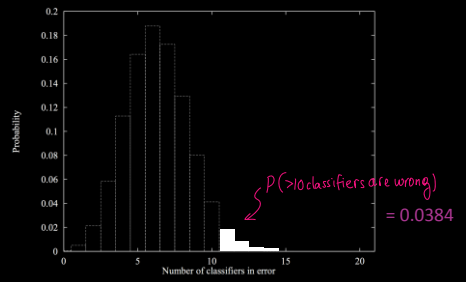with better performance.

## Teaser



## Teaser

Why do Ensembles work?
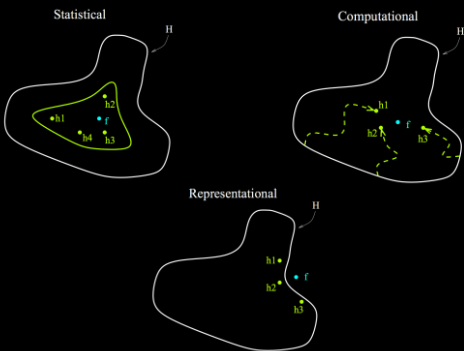
Hypothetical Classifier with $P_{error}=0.3$

Hypothetical Classifier with $P_{error}=0.3$



Hypothetical Classifier with $P_{error}=0.3$



$\sum P(>10 \text{ classifiers are wrong})$

$= 0.0384$

Why do Ensembles work?



BAGGING

## Bagging

Bagging (Boostrap aggregating).          (Breiman, 1996)

$\text{BAGGING}(S=((x_1,y_1),\dots,(x_m,y_m)))$
1  for $t \leftarrow 1$ to $T$ do
2      $S_t \leftarrow \text{BOOTSTRAP}(S) \triangleright$ i.i.d. sampling with replacement from $S$.
3      $h_t \leftarrow \text{TRAINCLASSIFIER}(S_t)$
4  return $h_S = x \mapsto \text{MAJORITYVOTE}((h_1(x),\dots,h_T(x)))$

## Bagging

Ensemble :

$$h_S(x) = \text{sign}\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right)$$

Bagging : Special case where we fix:

$$\alpha_t = 1 \qquad \text{and} \qquad h_t = \mathbb{L}(S_t)^*$$

* $\mathbb{L}$ is some learning algorithm

$S_t$ is a training set drawn from distribution $P(<x,y>)$
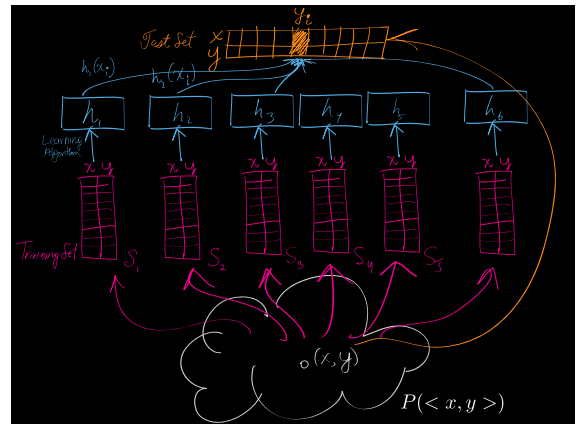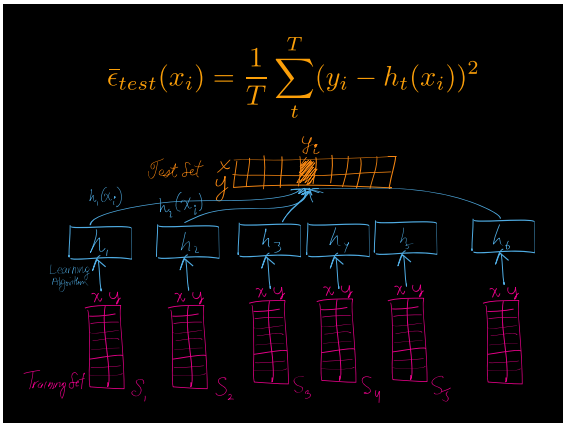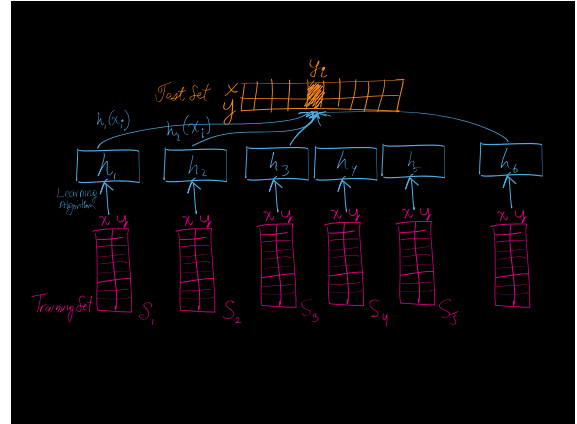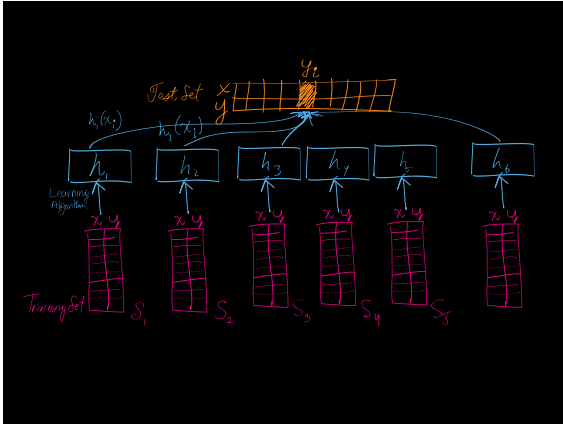
## Bias-Variance Tradeoff

## Generalization Error

Classification :

$$\epsilon_{test} = \frac{1}{n}\sum_{i}^{n} \text{Zero-One-Loss}(y_i, h(x_i))$$

Regression :

$$\epsilon_{test} = \frac{1}{n}\sum_{i}^{n} (y_i - h(x_i))^2$$

$$\bar{\epsilon}_{test}(x_i) = \frac{1}{T} \sum_t^T (y_i - h_t(x_i))^2$$



$$\bar{\epsilon}_{test}(x_i) = \frac{1}{T} \sum_t^T (y_i - h_t(x_i))^2$$

OR, as an expectation:

$$\mathbb{E}_S \left[ (y_i - h_S(x_i))^2 \right]$$

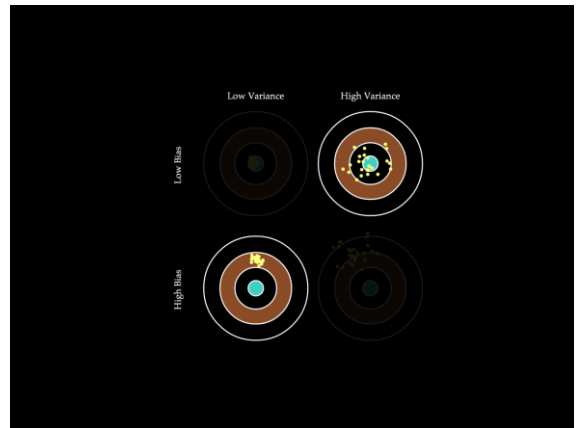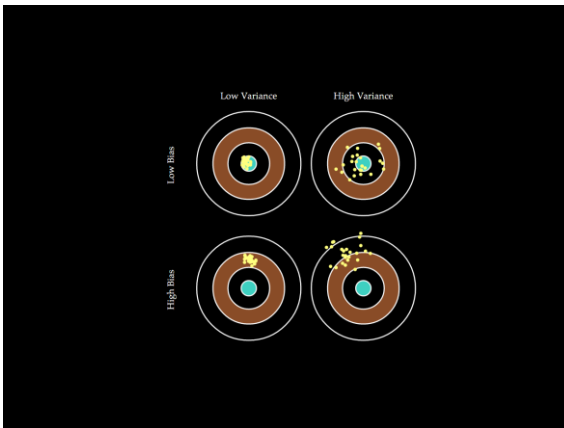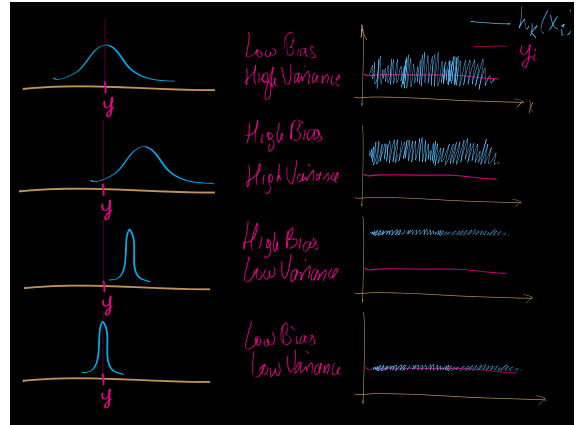For the entire test set:

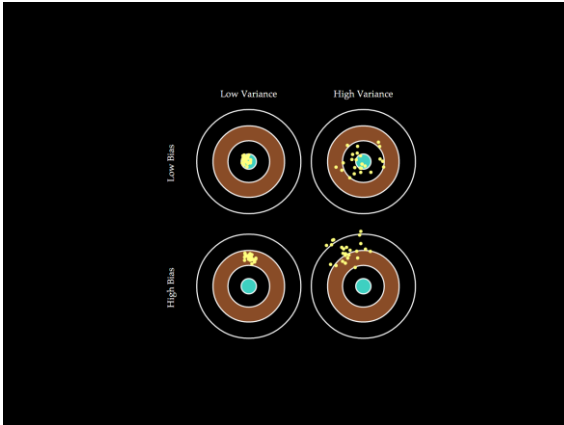$$\mathbb{E}_{X,Y} \mathbb{E}_S \left[ (y_i - h_S(x_i))^2 \right]$$

CLAIM:

$$\mathbb{E}_S \left[ (y_i - h_S(x_i))^2 \right] =$$

*bias²* $\qquad (y_i - \mathbb{E}_S[h_S(x_i)])^2 \ +$

*variance* $\qquad + \ \mathbb{E}_S[(h_s(x_i) - \mathbb{E}_S[(h_S(x_i))])^2]$

$$y = f(x)$$

$$y \in \{+1, -1\}$$

$$S_1$$

$$S_2$$

$$S_3$$

$$S_4$$

Consider: $h(x) = f(x)$

$y \in \{+1, -1\}$

Consider: $h(x) = f(x)$

$y \in \{+1, -1\}$

Consider: $h(x) = \mathbf{w}x + b$

$y \in \{+1, -1\}$

Consider: $h(x) = \mathbf{w}x + b$

$y \in \{+1, -1\}$

Consider: $h(x) = \mathbf{w}x + b$

$y \in \{+1, -1\}$

$$\mathbb{E}_S \left[ (y_i - h_S(x_i))^2 \right] =$$

$$(y_i - \mathbb{E}_S[h_S(x_i)])^2 \; +$$

$$+ \; \mathbb{E}_S[(h_s(x_i) - \mathbb{E}_S[(h_S(x_i))])^2]$$

## Label Noise

Noise-free:
$$y_i = f(x_i)$$

Regression:
$$y_i = f(x_i) + noise$$
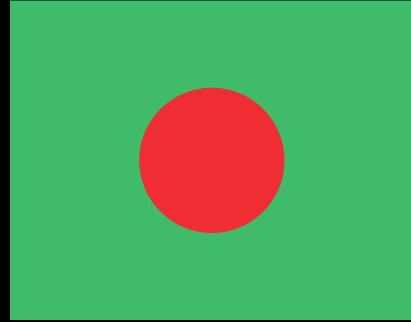$$y_i = f(x_i) + \mathcal{N}(0, \sigma^2)$$

Classification:
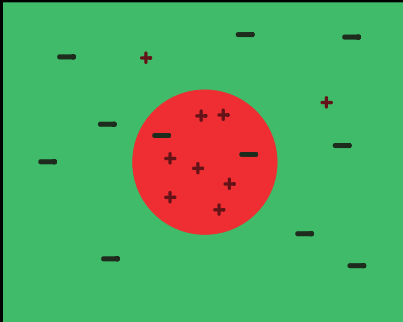$$y_i = noisy(f(x_i))$$

( *noisy()* switches label with probability *p* )

---

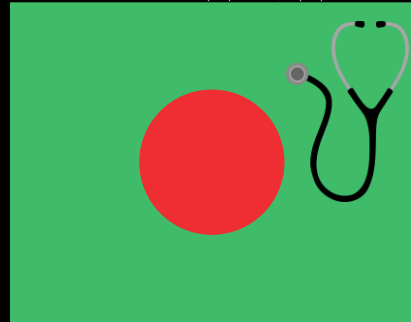$$y = \text{noisy}(f(x)) \quad \text{( flip sign with probability 0.25)}$$

$$y \in \{+1, -1\}$$

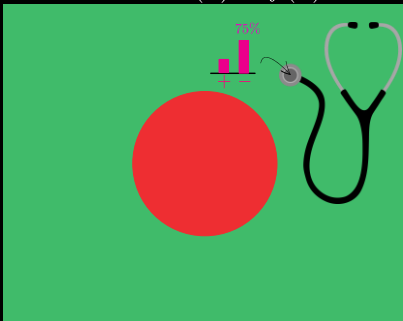---

$$S_1$$

$$y \in \{+1, -1\}$$

---

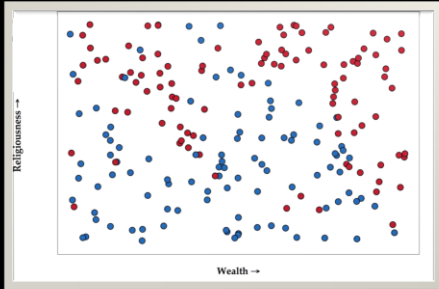Consider: $$h(x) = f(x) \quad y = \text{noisy}(f(x))$$

$$y \in \{+1, -1\}$$

---

Consider: $$h(x) = f(x) \quad y = \text{noisy}(f(x))$$

75%
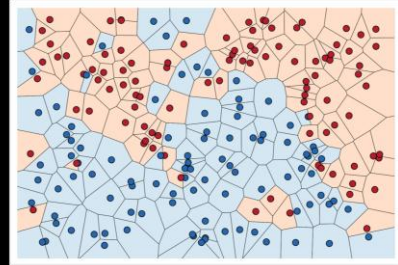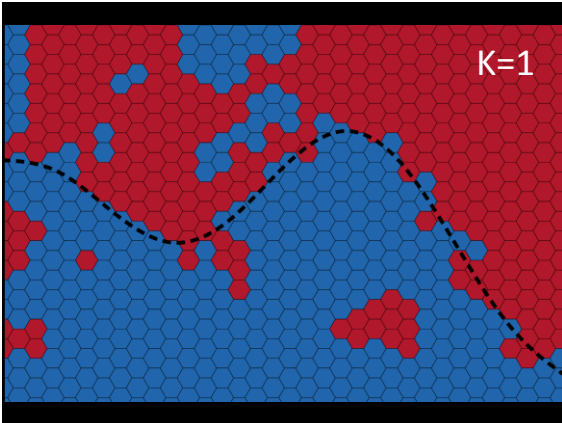
$$y \in \{+1, -1\}$$

---

## Example
*( kNN )*

Democrat vs Republican party association
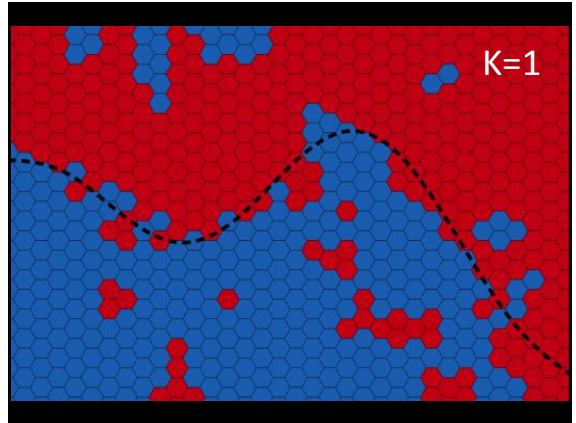


K=1



K=1



K=1



K=1



K=1

CLAIM:

$$\mathbb{E}_S\left[(y_i - h_S(x_i))^2\right] =$$

*bias²* $\qquad (y_i - \mathbb{E}_S[h_S(x_i)])^2 \; +$

*variance* $\quad + \mathbb{E}_S[(h_s(x_i) - \mathbb{E}_S[(h_S(x_i))])^2]$

USEFUL LEMMA:

$$\mathbb{E}[(\alpha - \mathbb{E}[\alpha])^2] = \mathbb{E}[\alpha^2] + \mathbb{E}[\alpha]^2$$

$$y_i = f(x_i) + \mathcal{N}(0, \sigma^2)$$

$$\mathbb{E}_S\left[(y_i - h_S(x_i))^2\right] =$$

*bias²* $\qquad (y_i - \mathbb{E}_S[h_S(x_i)])^2 \; +$

*variance* $\quad + \mathbb{E}_S[(h_s(x_i) - \mathbb{E}_S[(h_S(x_i))])^2]$

$$y_i = f(x_i) + \mathcal{N}(0, \sigma^2)$$

$$\mathbb{E}_S\left[(y_i - h_S(x_i))^2\right] =$$

*bias²*     $(f(x_i) - \mathbb{E}_S[h_S(x_i)])^2 +$

*variance*     $+ \mathbb{E}_S[(h_s(x_i) - \mathbb{E}_S[(h_S(x_i))])^2]$

*noise*     $+ \mathbb{E}_S[(f(x_i) - y_i)^2]$

---

$$y_i = f(x_i) + \mathcal{N}(0, \sigma^2)$$

$$\mathbb{E}_S\left[(y_i - h_S(x_i))^2\right] =$$

*bias²*     $(f(x_i) - \mathbb{E}_S[h_S(x_i)])^2 +$

*variance*     $+ \mathbb{E}_S[(h_s(x_i) - \mathbb{E}_S[(h_S(x_i))])^2]$

*noise*     $+ \sigma^2$

---

$$\mathbb{E}_S\left[(y_i - h_S(x_i))^2\right] =$$

*bias²*     $(y_i - \mathbb{E}_S[h_S(x_i)])^2 +$

*variance*     $+ \mathbb{E}_S[(h_s(x_i) - \mathbb{E}_S[(h_S(x_i))])^2]$

---

# BAGGING
## *revisited*



---

## Bagging

Bagging (Boostrap aggregating).                    (Breiman, 1996)

$\text{BAGGING}(S = ((x_1, y_1), \ldots, (x_m, y_m)))$
1    **for** $t \leftarrow 1$ **to** $T$ **do**
2        $S_t \leftarrow \text{BOOTSTRAP}(S)$ ▷ i.i.d. sampling with replacement from $S$.
3        $h_t \leftarrow \text{TRAINCLASSIFIER}(S_t)$
4    **return** $h_S = x \mapsto \text{MAJORITYVOTE}((h_1(x), \ldots, h_T(x)))$

# Why does it work?

---

## Bagging

Ensemble :

$$h_S(x) = \text{sign}\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right)$$

Bagging : Special case where we fix:

$$\alpha_t = 1 \quad \text{and} \quad h_t = \mathbb{L}(S_t)^*$$

$^*$ $\mathbb{L}$ is some learning algorithm

$S_t$ is a training set drawn from distribution     $P(<x, y>)$

## Bagging

Bagging Ensemble :

$$h_S(x) = \text{sign}\left(\sum_{t=1}^{T} h_t(x)\right)$$

What happens to *bias* and *variance?*

## Bagging

Bagging Ensemble ( regression ) :

$$h_S(x) = \frac{1}{T}\sum_{t=1}^{T} h_t(x)$$

*bias²*      $(y_i - \mathbb{E}_S[h_S(x_i)])^2$

*variance*    $\mathbb{E}_S[(h_s(x_i) - \mathbb{E}_S[(h_S(x_i))])^2]$
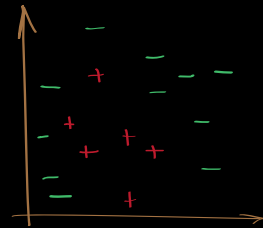
## Bagging

What happens to *bias* and *variance?*

$$\text{Bias}(h_s, x_i) =$$

$$\text{Var}(h_s, x_i) \approx$$

Bagging has approximately the same bias, but reduces variance of individual classifiers!

## Bagging



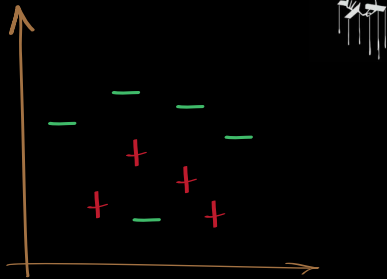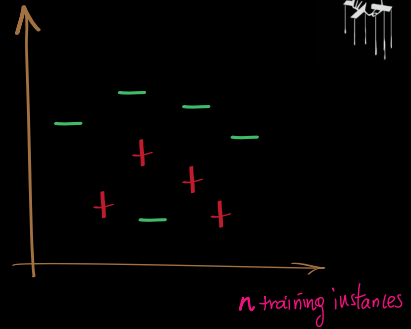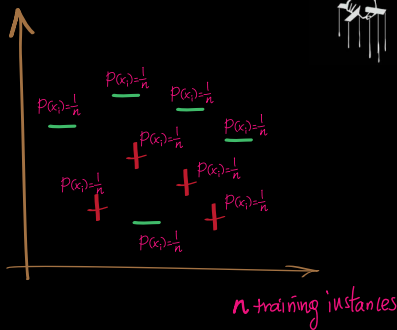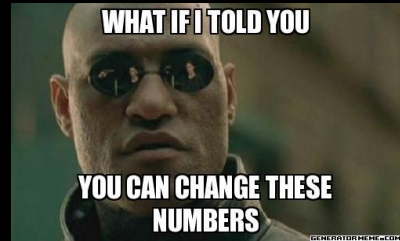## Bagging as a "Training set manipulator"

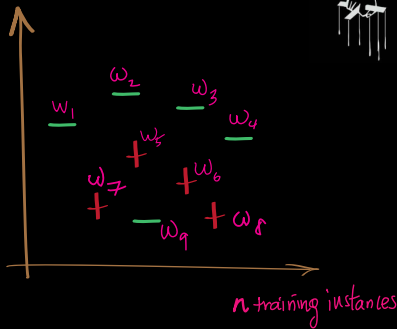## Bagging as a "Training set manipulator"

Bagging as a "Training set manipulator"



Bagging as a "Training set manipulator"

$n$ training instances



Bagging as a "Training set manipulator"

$p(x_i) = \frac{1}{n}$

$n$ training instances



Bagging as a "Training set manipulator"



WHAT IF I TOLD YOU

YOU CAN CHANGE THESE NUMBERS

Bagging as a "Training set manipulator"

$w_1$, $w_2$, $w_3$, $w_4$, $w_5$, $w_6$, $w_7$, $w_8$, $w_9$

$n$ training instances



## Ensemble

Problem : given *T* binary classification hypotheses ($h_1$,…, $h_T$), **find** a combined classifier:

$$h_S(x) = \text{sign}\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right)$$

with better performance.

Teaser

BOOSTING